

Dear Sir:

Transmitted herewith for filing is the patent application of

Inventor(s): D.A. Burton; R.L. Morton; and E. Webman  
For: **METHOD, SYSTEM, AND PROGRAM FOR SELECTING ONE OF MULTIPLE PATHS TO COMMUNICATE WITH A DEVICE**

Enclosed are:

- ☒ 5 No. of Sheets of Drawings Sheet(s) of drawings ( ☒ informal) + 0 extra copies;
- ☒ 26 pages of Application; 15 pages of specification, 1 page of abstract
- ☐ An assignment of the invention to International Business Machines Corporation. ( ☐ Will follow.)
- ☐ An associate power of attorney.
- ☐ A verified statement to establish small entity status under 37 CFR 1.9 and 1.27.
- ☒ Declaration and Power of Attorney. ( ☐ Will follow.)
- ☐ Certified copy of Patent Application No. filed from which priority is claimed under 35 U.S.C. §119.
- ☐ IDS enclosed. ☐ with references.

CALCULATION OF FEES							
ITEM		NO. OF CLAIMS FILED MINUS BASE*		NO. OF CLAIMS OVER BASE	X SM/LG ENTITY FEE	\$ AMOUNT	\$ FEE
A	TOTAL CLAIMS FEE	42	- 20* =	22	X \$9 or \$18	\$396	
B	INDEPENDENT CLAIMS FEE**	3	- 3* =	0	X \$39 or \$78	\$0	
C	SUBTOTAL - ADDITIONAL CLAIMS FEE (ADD FINAL COLUMN IN LINES A + B)						396
D	MULTIPLE-DEPENDENT CLAIMS FEE						

Please charge Dep. Acct. No. 50-0585 in the amount of \$ **A copy of this sheet is enclosed.**

- ☒ A check in the amount of \$ 1,086 to cover the filing fee is enclosed.
- ☒ Check for \$ 40 covering the Recordation of Assignment fee enclosed.
- ☒ The Commissioner is hereby authorized to charge payment of the following fees associated with this communication or credit any overpayment to Deposit Account No. 50-0585 . **A copy of this sheet is enclosed.**
  - ☒ Any additional filing fees required under 37 CFR 1.16.
  - ☒ Any patent application processing fees under 37 CFR 1.17.
- ☐ The Commissioner is hereby authorized to charge payment of the following fees during the pendency of this application or credit any overpayment to Deposit Account No. 50-0585. **A copy of this sheet is enclosed.**
  - ☐ Any patent application processing fees under 37 CFR 1.17.
  - ☐ The issue fee set in 37 CFR 1.18 at or before mailing of the Notice of Allowance, pursuant to 37 CFR 1.311(b).
  - ☐ Any filing fees under 37 CFR 1.16 for presentation of extra claims.



Respectfully submitted,

David W. Victor  
Registration No. 39,867

Direct All Correspondence to:  
David W. Victor  
KONRAD RAYNES & VICTOR LLP  
1180 S. Beverly Drive; Suite 501  
Los Angeles, CA 90035

Direct Telephone Calls to:  
(310) 556-7983



METHOD, SYSTEM, AND PROGRAM FOR SELECTING ONE OF  
MULTIPLE PATHS TO COMMUNICATE WITH A DEVICE

Cross-Reference to Related Applications

5        This application is related to the following co-pending and commonly-assigned patent applications, all of which are filed on the same date herewith, and all of which are incorporated herein by reference in their entirety:

10        “Method, System, And Program For Determining A Number of Write Operations to Execute”, to David A. Burton, Robert L. Morton, and Erez Webman, having attorney docket no. TUC9-2000-0015US1, and  
      “Method, System, And Program For Remote Copy in an Open Systems Environment” to David A. Burton, Robert L. Morton, and Erez Webman, having attorney docket no. TUC9-2000-0016US1.

15                                    BACKGROUND OF THE INVENTION

1.        Field of the Invention

      The present invention relates to a system, method, and program for selecting a path to use to communicate to a device to improve transmission performance.

20    2.        Description of the Related Art

      Two systems communicating over a network may each include multiple ports, thus providing multiple paths across which data can be communicated. In certain prior art systems, a path may be selected according to a round robin or other predefined path rotation technique or a single default path is used for all operations. However, such  
25 techniques do not attempt to optimize path selection when there are multiple available paths.

	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100	2101	2102	2103	2104	2105	2106	2107	2108	2109	2110	2111	2112	2113	2114	2115	2116	2117	2118	2119	2120	2121	2122	2123	2124	2125	2126	2127	2128	2129	2130	2131	2132	2133	2134	2135	2136	2137	2138	2139	2140	2141	2142	2143	2144	2145	2146	2147	2148	2149	2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	2160	2161	2162	2163	2164	2165	2166	2167	2168	2169	2170	2171	2172	2173	2174	2175	2176	2177	2178	2179	2180	2181	2182	2183	2184	2185	2186	2187	2188	2189	2190	2191	2192	2193	2194	2195	2196	2197	2198	2199	2200	2201	2202	2203	2204	2205	2206	2207	2208	2209	2210	2211	2212	2213	2214	2215	2216	2217	2218	2219	2220	2221	2222	2223	2224	2225	2226	2227	2228	2229	2230	2231	2232	2233	2234	2235	2236	2237	2238	2239	2240	2241	2242	2243	2244	2245	2246	2247	2248	2249	2250	2251	2252	2253	2254	2255	2256	2257	2258	2259	2260	2261	2262	2263	2264	2265	2266	2267	2268	2269	2270	2271	2272	2273	2274	2275	2276	2277	2278	2279	2280	2281	2282	2283	2284	2285	2286	2287	2288	2289	2290	2291	2292	2293	2294	2295	2296	2297	2298	2299	2300	2301	2302	2303	2304	2305	2306	2307	2308	2309	2310	2311	2312	2313	2314	2315	2316	2317	2318	2319	2320	2321	2322	2323	2324	2325	2326	2327	2328	2329	2330	2331	2332	2333	2334	2335	2336	2337	2338	2339	2340	2341	2342	2343	2344	2345	2346	2347	2348	2349	2350	2351	2352	2353	2354	2355	2356	2357	2358	2359	2360	2361	2362	2363	2364	2365	2366	2367	2368	2369	2370	2371	2372	2373	2374	2375	2376	2377	2378	2379	2380	2381	2382	2383	2384	2385	2386	2387	2388	2389	2390	2391	2392	2393	2394	2395	2396	2397	2398	2399	2400	2401	2402	2403	2404	2405	2406	2407	2408	2409	2410	2411	2412	2413	2414	2415	2416	2417	2418	2419	2420	2421	2422	2423	2424	2425	2426	2427	2428	2429	2430	2431	2432	2433	2434	2435	2436	2437	2438	2439	2440	2441	2442	2
--	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	---

## 5

10

15

15

20

25

disabled and removed from the potential selection pool due to relatively poor performance for the path. In this way, preferred embodiments provide an algorithm for optimizing path performance.

5                                    BRIEF DESCRIPTION OF THE DRAWINGS

Referring now to the drawings in which like reference numbers represent corresponding parts throughout:

FIG. 1 is a block diagram illustrating a computing environment in which preferred embodiments are implemented; and

10            FIGs. 2a, b illustrate an example of data structures used in accordance with the preferred embodiments of the present invention; and

FIGs. 3, 4, and 5 illustrate logic implemented in a controller to select one of multiple paths to use to transfer data to a remote controller in accordance with the preferred embodiments of the present invention.

15                                    DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

In the following description, reference is made to the accompanying drawings which form a part hereof and which illustrate several embodiments of the present invention. It is understood that other embodiments may be utilized and structural and  
20    operational changes may be made without departing from the scope of the present invention.

FIG. 1 illustrates a computing environment in which preferred embodiments are implemented. Hosts 4a, b may comprise any computing device known in the art, including servers through which other client computers can access storage or clients. The  
25    hosts 4a, b each include at least one adaptor, such as a Fibre Channel or Small Computer System Interface (SCSI) adaptor card or any other network adaptor card known in the art. The host adaptors allow the hosts 4a, b to communicate with storage controllers 6a, b via switches 8a, b. The switches 8a, b may comprise the International Business Machines

Corporation (IBM) Fibre Channel Storage Hub or Switch, the IBM SAN Fibre Channel Switch, or any other switching device known in the art. Each switch 8a, b has a port connecting to a network 12, which may comprise any local area network, wide area network, the Internet or any other network system. The network 12 may use routers and  
5 switches to dynamically determine the data path through the network 12.

In the described implementations, a primary controller 6a includes interface cards 14a and b having ports 16a, b, c, d and a secondary controller 6b includes interface cards 18a and b having ports 20a, b, c, d. Primary controller 6a would communicate with the secondary controller 6b via one of the ports 16a, b, c, d, switch 8a, the network 12, switch  
10 8b, and then one of the ports 20a, b, c, d on the secondary controller 6b. Thus, the primary controller 6a can select one of sixteen paths to communicate with the secondary controller 6b, i.e., one of the ports 16a, b, c, d paired with one of the ports 20a, b, c, d. In alternative embodiments, each of the controllers 6a, b may include a different number of interface cards having a different number of ports to provide more or less communication  
15 paths therebetween.

In the preferred embodiments, the secondary storage 10b maintains a mirror copy of specified data volumes in the primary storage 10a. During an establishment phase, a relationship is established between primary volumes in the primary storage 10a and corresponding secondary volumes in the secondary storage 10b that mirror the primary  
20 volumes. After this relationship is established, the primary controller 6a will write any updates from hosts 4a, b to primary volumes to the secondary controller 6b to write to the secondary volumes in the secondary storage 10b .

The primary and secondary controllers 6a, b may include IBM Peer-to-Peer Remote Copy (PPRC), Extended Remote Copy (XRC) software, or other vender  
25 shadowing software to allow communication between the controllers 6a, b to coordinate data shadowing. In such embodiments, the controllers 6a, b may comprise large scale storage controllers, such as the IBM 3990 and Enterprise Storage System class controllers.\*\* In open system embodiments, the primary and secondary controllers 6a, b

may comprise controllers from different vendors of different models, etc., and may not include any specialized protocol software for performing the backup operations. Further, the controllers may include any operating system known in the art, including the Microsoft Corporation Windows and NT operating systems.\*\* In open systems

5   embodiments, the primary controller 6a can use commonly used write commands, such as SCSI write commands, to copy the primary volumes to the secondary volumes in the secondary storage 10b. In such open system embodiments, the secondary controller 6b does not need special purpose software to coordinate the shadowing activities with the primary controller 6b as the primary controller 6a accomplishes the shadowing by using

10   standard write commands. Further, in such open systems, the primary and secondary controllers 6a, b may comprise any controller device known in the art and the primary and secondary controllers 6a, b may be of different models and model types, and even of different classes of storage controllers.

Because there are multiple paths through which the primary controller 6a may

15   communicate with the secondary controller 6b over a network 12, preferred embodiments provide an algorithm for the primary controller 6a to use when selecting a path from one of the ports 16a, b, c, d in the primary controller 6a to one of the ports 20a, b, c, d in the secondary controller 6b. The primary controller 6a may use this path selection algorithm when determining a path to use to communicate updates from a host 4a to primary

20   volumes to be shadowed in secondary volumes in the secondary storage 10b.

Below are data structures, that are used by the path selection algorithm shown in FIGs. 3, 4, and 5:

Cumulative Transfer Time (cumulativeXferTime): An array data structure, as shown in FIG. 2a, has an entry for each of the sixteen paths and each of the block

25   size ranges. Each of the sixteen rows corresponds to one port 16a, b, c, d paired with one port 20a, b, c, d. Each of the columns provides a block size range for the size of the update transferred down the path. Thus, the cumulativeXferTime array provides the total transfer time for all transfers in a measurement period down one

of the sixteen paths within one of the three block size ranges, e.g., less than nine blocks, between nine and sixty-four blocks and more than sixty-four blocks.

Number Transfers (numXfers): An array data structure having an entry for one of the sixteen paths and one of the block size ranges, i.e., the same number of entries and column and row labels as the cumulative transfer time array. The number transfers array accumulates the number of transfers in a measurement period down one of the sixteen paths having one of the three block size ranges. If a path for a particular block size range is disabled, then the entry in the numXfers array for the path and block size range will maintain a counter indicating the number of transfers for which the path will remain disabled for that block size range. Once this counter is decremented to zero, the path for the block size range will be enabled and available for use.

Transfer Count (xferCount): A one column array, as shown in FIG. 2b, having one entry for each block size range. Each entry provides the number of transfers across all paths for a given block size during the measurement period. A measurement period ends for a given block size when the value in one of the entries in the xferCount array reaches a predetermined value, such as 128 or any other selected number.

Path Enabled (pathEnabled): An array data structure having an entry for one of the sixteen paths and one of the block size ranges, i.e., the same number of entries and column and row labels as the cumulative transfer time array. Each entry is a boolean value indicating whether the preferred path is enabled, i.e., whether the controller 6a can select this path to use for a write operation to the secondary controller 6b having an update within the block size range for the entry.

The preferred path algorithm shown in FIGs. 3, 4, and 5 uses the above data structures to gather performance information for each path and for updates within one of three different block size ranges. In preferred embodiments, data is gathered for a



measurement period of  $n$  transfers, which in the described implementation is set to 128. After  $n$  transfers for one of the block sizes, the preferred path algorithm analyzes the data to determine whether to disable or enable the paths for that block size. If a given path is disabled for a block size, then the primary controller 6a would not select such disabled  
5 path for an update within the block size range and would only select one of the enabled paths for the block size range. The path selection algorithm may use any selection technique known in the art, e.g., round robin, etc., for selecting from one of the enabled paths to use to transfer an update within a block size range.

FIGs. 3, 4, and 5 illustrate path selection logic implemented as code and executed  
10 by the primary controller 6a. With respect to FIG. 3, control begins at block 100 with the primary controller 6a initiating shadowing operations to write any updates to primary volumes in the primary storage 10a to the secondary controller 6b after the establishment of a relationship between primary volumes in the primary storage 10a and secondary volumes in the secondary storage 10b. The primary controller 6a then initializes all the  
15 arrays (at blocks 102 and 104) by setting all the entries in the cumulative transfer time (cumXferTime) array, number transfers array (numXfers), and the transfer count array (xferCount) to zero. All entries in the path enabled array (pathEnabled) are set to "on", indicating that initially all paths for all block size ranges are available for selection.

After initializing all the arrays, the primary controller 6a would wait (at block  
20 130) for any write operation  $i$  comprising an update to a primary volume from hosts 4a, b. In response, the primary controller 6a would determine (at block 132) the block size range  $k$  (in the described implementation there are three possible ranges) including the size of the update in write  $i$ . In SCSI embodiments, the primary controller 6a can determine the size of the update from the transfer length field in the write command  
25 descriptor block (CDB). The primary controller 6a would then select (at block 134) from the path enabled array (pathEnabled) an enabled path, i.e., an entry for a path that has an "on" value, for the determined block size range  $k$ . The primary controller 6a may use any selection procedure known in the art for selecting one enabled path, round robin, etc. The

primary controller 6a would then start (at block 136) timer  $i$  for write  $i$  and send (at block 138) the write  $i$  to the secondary controller 6b to apply to the secondary volume in secondary storage 10b. The update would also be applied to the primary volume in primary storage 10a.

5           At block 170, the primary controller 6a waits for a response to one of the outstanding write operations, or write  $i$  which was initially sent down path  $j$ . The primary controller 6a may maintain information associating a selected path  $j$ , a write  $i$ , and the timer  $i$ . When the response for a write is received, the primary controller 6a can then use this information to determine the write  $i$  for which the response is received, and also the  
10 timer  $i$  and path  $j$  for write  $i$ . Upon receiving a response from the secondary controller 6b that the write  $i$  completed, the primary controller 6a would stop (at block 172) timer  $i$  for the completed write  $i$ . If write  $i$  did not successfully complete (at block 174), then an error mode would occur (at block 176). Otherwise, if the write  $i$  was successful, then the primary controller 6a would (at block 178) add the time value in timer  $i$  to the entry in the  
15 cumulative transfer time array (cumXferTime) for path  $j$  and block size range  $k$ , or entry  $j, k$ , where the update of write  $i$  falls within block size range  $k$ . The entry  $j, k$  in the number of transfers array (numXfers), indicating the number of writes completed down path  $j$  within the block size range  $k$ , is incremented (at block 180) by one. Further the  $k$ th entry in the transfer count array (xferCount) is incremented (at block 182) by one, indicating  
20 the number of writes completed for that particular block size range  $k$ . If, at block 184, the value in the incremented entry  $k$  in the transfer count array (xferCount) is equal to 128, or any other measurement period integer value, then control proceeds to block 200 in FIG. 4; otherwise, the logic ends.

25           The preferred logic of FIG. 3 accumulates the time for write operations for each path by block size range so that any comparisons of time for the paths will primarily be based on network transmission factors. Because the time is accumulated for a same transfer size range, the size effect of the transfer on the performance of the secondary controller 6b in affecting the secondary controller 6b load will be relatively constant.

Thus, the primary difference in the performance among paths within a same block size range will be network transmission delays associated with the paths as the load the write operation imposes on the secondary controller 6b remains relatively constant for the cumulative time measurements within a block size range.

5           The preferred logic of FIGs. 4 and 5 is implemented if one of the entries for a block size in the transfer count (xferCount) array equals the measurement period, which in the described implementation is 128. Thus, a path enablement decision is made for the paths in a particular block size range  $k$  after the data points collected for a range  $k$  reaches the measurement period. The measurement period should be selected sufficiently large to  
10 allow useful statistical information to be gathered for the paths in a particular block size range, yet not too large such that enablement adjustments would be unduly delayed. Lengthy delays in adjustments to the path enablement settings would allow the continued selection of poor performing paths that would otherwise be disabled during an adjustment.

15           If the measurement period, e.g., 128 write operations, has been reached for the block size range  $k$  entry in the transfer count array (xferCount), then the primary controller 6a, at blocks 200-214 in FIG. 4, will decrement the number transfers counter (numXfer) for each disabled path and block size range. As discussed, if a path is disabled for a block size range, then the numXfer entry for that path and block size maintains a  
20 counter that is decremented for each transfer. Once the entry is decremented to zero, the path is reenabled and made available for use. For each path  $m$  where  $m$  equals 1 to 16 (or any other alternative number of possible paths between the devices), the primary controller 6a will perform steps 204 to 212 within the loop to adjust the number transfer counter (numXfer) for disabled paths for the given block size range  $k$ . At block 204, the  
25 primary controller 6a determines whether entry  $m, k$  in the path enabled array (pathEnabled) is "off". If so, then the primary controller 6a subtracts (at block 206) 128 from the entry  $m, k$  in the number of transfers array (numXfers). In the described implementation, the primary controller 6a subtracts a number of completed transfers

every 128 transfers for the block size range, thus saving processing cycles by not having to perform the logic of FIG. 4 after every write operation.

If (at block 210) the entry  $m, k$  in the number transfers array (numXfers) was just decremented to zero at block 206, then the primary controller 6a sets (at block 212) the entry  $m, k$  in the path enabled array (pathEnabled) to "on" making that path  $m$  for block size range  $k$  available for use. If the entry  $m, k$  is already enabled (no branch of block 204) or the counter is not decremented to zero (no branch of 210) or after the entry  $m, k$  is enabled (at block 212), then control transfers to block 214 where the primary controller 6a will proceed back to block 200 if there are further paths to consider. Otherwise, control proceeds to block 250 in FIG. 5.

At blocks 250-266 in FIG. 5, the primary controller 6a determines whether to disable any paths for the block size range  $k$  whose entry in the transfer count array (xferCount) reached the measurement period value, e.g., 128. At block 250, the primary controller determines an average transfer time  $m, k$  for each path  $m$  for the block size range  $k$  by dividing the entry  $m, k$  for path  $m$  and block size range  $k$  in the cumulative transfer time array (cumXferTime), indicating the total transfer time for all completed writes for path  $m$  and block size range  $k$ , by the number of transfers for that path  $m$  and block size range  $k$  that resulted in the total transfer time indicated in entry  $m, k$  in cumXferTime. From the calculated average transfer times calculated at block 250, the best average total transfer time for block size range  $k$  is then determined (at block 252). The primary controller 6a then performs a loop between blocks 254 and 266 to determine whether to disable each path  $m$  for the block size range  $k$ .

If (at block 256) the average transfer time  $m, k$  for path  $m$  is not 15% longer than the best average transfer time, then the path  $m$  for block size range  $k$  is not disabled and control proceeds (at block 266) to consider the next  $(m + 1)$ th path in block size range  $k$ . Otherwise, if the transfer time of path  $m$  is 15% longer than the best average transfer time, then the path  $m$  for block size range  $k$  is disabled (at block 258) by setting the entry  $m, k$  in the path enabled array (pathEnabled) to "off". If the average transfer time  $m, k$  for

path  $m$  is between 15% and 25% longer than the best average transfer time (at block 260), then the entry  $m, k$  for path  $m$  in the number of transfers array (numXfers) is set (at block 262) to 4096, indicating the number of transfers for block size range  $k$  before the path  $m$  for block size rang  $k$  is reset to enabled and available for use. Otherwise, if the average  
5 transfer time  $m, k$  for path  $m$  is more than 25% longer than the best average transfer time, then the entry  $m, k$  in numXfers is set (at block 264) to 8192. In preferred embodiments, the number of transfers for which a path is disabled is a multiple of the measurement period, to ensure that subtraction of the measurement period from the disablement counter will eventually produce zero.

10 With the logic of FIG. 5, the path remains disabled for a longer number of transfers if the performance for the path is 25% worse than the best average performance, i.e., the average transfer time for a path  $m$  and given block size range  $k$  takes more than 25% longer than the best average transfer time as opposed to if the average transfer time for path  $m$  is only 15% to 25% worse than the best average transfer time. In other words,  
15 the more degraded the performance for a path  $m$  and block size range  $k$ , the longer the path will be disabled and taken off-line.

The preferred embodiment algorithm for selecting paths optimizes overall performance when multiple paths are available between the two devices by removing those paths from selection whose performance is appreciably worse than that of other  
20 paths. Those paths removed from selection are not available for selection in the round robin path selection process to use for a transfer operation between the devices. Preferred embodiments gather performance data and periodically analyze the data to determine whether to adjust the enablement or disablement setting for each path for a given block size range. Certain of the performance problems associated with a path may only be  
25 temporary. For this reason, it is desirable to occasionally enable previously disabled paths so that transient conditions do not permanently remove a path from the selection process. In this way, temporary bottlenecks in the system are avoided and the best pair of source and destination ports 16a, b, c, d, and 20a, b, c, d (FIG. 1) are selected for

communication between the primary 6a and secondary 6b controllers. In further embodiments, if there is no enabled path available for a particular transfer operation, then one of the disabled paths may be selected and used.

In open systems embodiments, the primary controller 6a is able to determine the path performance to the secondary controller 6b without having to establish a special communication protocol, which would require software on both the primary 6a and secondary 6b controllers, e.g., IBM PPRC and XRC software. Instead, in preferred embodiments, the primary controller 6a may write updates to the secondary controller 6b using standard SCSI commands, and can select an optimal path based on acknowledgment information the secondary controller 6b returns under the SCSI protocol. In this way, the secondary controller 6b does not have to know that it is being monitored as the primary controller 6a independently handles the monitoring. In alternative embodiments where the primary 6a and secondary 6b controllers include specialized shadowing software, there may be additional communications to perform path selection optimization.

### Conclusion

The following describes some alternative embodiments for accomplishing the present invention.

The preferred embodiments may be implemented as a method, apparatus or program using standard programming and/or engineering techniques to produce software, firmware, hardware, or any combination thereof. The programs defining the functions of the preferred embodiment can be delivered to a computer via a variety of information bearing media, which include, but are not limited to, computer-readable devices, programmable logic, memory devices (e.g., EEPROMs, ROMs, PROMs, RAMs, SRAMs, etc.) carriers, or media, such as a magnetic storage media, "floppy disk," CD-ROM, a file server providing access to the programs via a network transmission line, wireless transmission media, signals propagating through space, radio waves, infrared

signals, etc. Of course, those skilled in the art will recognize that many modifications may be made to this configuration without departing from the scope of the present invention. Such signal-bearing media, when carrying computer-readable instructions that direct the functions of the present invention, represent alternative embodiments of the  
5 present invention.

In preferred embodiments, one primary controller 6a shadowed data on a single secondary controller 6b. However, in alternative arrangements, one primary controller may shadow data at multiple secondary controllers or multiple primary controllers may shadow data at one or more secondary controllers.

10 The preferred logic of FIGs. 3-5 describes specific operations occurring in a particular order. In alternative embodiments, certain of the logic operations may be performed in a different order, modified or removed and still implement preferred embodiments of the present invention. Moreover, steps may be added to the above described logic and still conform to the preferred embodiments. Further, operations  
15 described herein may occur sequentially or certain operations may be processed in parallel.

In preferred embodiments, the primary and secondary controllers communicated over a fibre channel interface using SCSI commands. In alternative embodiments, different command protocols may be utilized. For instance, the ESCON protocol may be  
20 used for the channel communications and the count-key-data (CKD) protocol may be used for the input/output (I/O) commands.

Preferred embodiments were described with respect to a storage system in which data from a primary site is shadowed at a secondary site as part of a data backup and recovery system. However, the preferred method, system, and program for selecting an  
25 optimal path between a primary 6a and secondary 6b controller for data shadowing operations may apply to any two devices having multiple paths therebetween. For instance, the preferred embodiment path selection algorithm may apply to any situation where one system is selecting from one of multiple paths to another system and, in

particular, to paths that are directed through a network that may provide additional routing and switching of the paths.

In preferred embodiments a write operation including an update was transmitted down the selected path. In alternative embodiments, any type of data may be  
5 communicated via the selected path.

Further, in preferred embodiments the paths traversed a network. In alternative embodiments the paths may comprise point-to-point paths communication lines between the two devices.

In preferred embodiments, specific values were specified for the measurement  
10 period (128) and number of writes to keep a path disabled (4096 or 8192). In alternative embodiments different values may be used. Further, there may be additional thresholds for providing additional transfers a path may remain disabled based on additional performance criteria.

In summary, preferred embodiments disclose a method, system, and program for  
15 selecting one of multiple data paths to a device. A selection is made of one of multiple paths indicated as enabled to transmit data. A path is indicated as enabled or disabled. Transfer time data is gathered for each enabled path capable of being selected. Paths having transfer time data satisfying a threshold are indicated as disabled. Paths indicated as disabled are not capable of being selected to use to transmit data.

20 The foregoing description of the preferred embodiments of the invention has been presented for the purposes of illustration and description. It is not intended to be exhaustive or to limit the invention to the precise form disclosed. Many modifications and variations are possible in light of the above teaching. It is intended that the scope of the invention be limited not by this detailed description, but rather by the claims  
25 appended hereto. The above specification, examples and data provide a complete



5

**\*\*Enterprise Storage Server and ESCON are registered trademarks and Fibre Channel Raid Storage Controller is a trademark of IBM; Windows and Windows NT are registered trademarks of Microsoft Corporation.**

WHAT IS CLAIMED IS:

1           1.       A method for selecting one of multiple data paths to a device, comprising:  
2           selecting one of multiple paths indicated as enabled to transmit data, wherein a  
3           path is indicated as enabled or disabled;  
4           gathering transfer time data for each enabled path capable of being selected; and  
5           indicating paths as disabled having transfer time data satisfying a threshold,  
6           wherein paths indicated as disabled are not capable of being selected to use to transmit  
7           data.

1           2.       The method of claim 1, further comprising:  
2           indicating one disabled path as enabled after performing a threshold number of  
3           transfer operations.

1           3.       The method of claim 2, further comprising:  
2           disabling the path for a first threshold number of transfer operations if the path has  
3           a transfer data time satisfying a first threshold; and  
4           disabling the path for a second threshold number of transfer operations if the path  
5           has a transfer data time satisfying a second threshold.

1           4.       The method of claim 1, wherein transfer time data is gathered by path and  
2           transfer size, wherein a path is disabled for a given transfer size and wherein one path  
3           disabled for one transfer size is capable of being enabled for at least one other transfer  
4           size.

1           5.       The method of claim 1, wherein gathering transfer time data further  
2           comprises:

[illegible]

1           9.       The method of claim 5, wherein the transfer time is measured from the  
2   time the transfer is sent to the device to the time a response is received from the device

3 indicating that the transfer completed, further comprising adding the transfer time for a  
4 transfer transmitted down the path to the cumulative transfer time for the path.

1 10. The method of claim 5, further comprising:  
2 for each enabled path, determining a best average transfer time from the average  
3 cumulative transfer times for all paths, wherein determining whether the average  
4 cumulative transfer time for one path satisfies the threshold comprises determining  
5 whether the average cumulative transfer time for the path exceeds the best average  
6 transfer time by a percentage amount.

1 11. The method of claim 10, wherein determining whether the average  
2 cumulative transfer time satisfies the threshold further comprises disabling the path for a  
3 first number of transfer operations if the average cumulative transfer time for the path  
4 exceeds the best average transfer time by a first percentage amount and disabling the path  
5 for a second number of transfer operations if the average cumulative transfer time for the  
6 path exceeds the best average transfer time by a second percentage amount.

1 12. The method of claim 1, wherein the multiple paths comprise multiple  
2 paths between a first controller and a second controller, and wherein one path is selected  
3 to transmit updates to a primary storage area managed by the first controller to the second  
4 controller to store in a secondary storage area.

1 13. The method of claim 11, wherein transfer time data is gathered by path  
2 and a size of the update, wherein a path is disabled for a given update size and wherein  
3 the path is capable of being enabled for at least one other update size.

1 14. The method of claim 1 wherein the paths extend through a network.

1           15.     A system for selecting one of multiple data paths to a device, comprising:  
2                 means for selecting one of multiple paths indicated as enabled to transmit data,  
3     wherein a path is indicated as enabled or disabled;  
4                 means for gathering transfer time data for each enabled path capable of being  
5     selected; and  
6                 means for indicating paths as disabled having transfer time data satisfying a  
7     threshold, wherein paths indicated as disabled are not capable of being selected to use to  
8     transmit data.

1           16.     The system of claim 15, further comprising:  
2                 means for indicating one disabled path as enabled after performing a threshold  
3     number of transfer operations.

1           17.     The system of claim 16, further comprising:  
2                 means for disabling the path for a first threshold number of transfer operations if  
3     the path has a transfer data time satisfying a first threshold; and  
4                 means for disabling the path for a second threshold number of transfer operations  
5     if the path has a transfer data time satisfying a second threshold.

1           18.     The system of claim 15, wherein transfer time data is gathered by path and  
2     transfer size, wherein a path is disabled for a given transfer size and wherein one path  
3     disabled for one transfer size is capable of being enabled for at least one other transfer  
4     size.

1           19.    The system of claim 15, wherein the means for gathering transfer time data  
2 further comprises:

3               means for gathering a cumulative transfer time for all transfer operations during a  
4 measurement period through the path and a cumulative number of the transfer operations  
5 during the measurement period for each enabled path; and

6               means for determining the average cumulative transfer time for the measurement  
7 period by dividing the cumulative time by the cumulative number of transfers for each  
8 enabled path, wherein the path is indicated as disabled if the average cumulative transfer  
9 time for the path satisfies the threshold.

1           20.    The system of claim 19, wherein the measurement period comprises a  
2 number of transfer operations for all paths, wherein the determination to disable paths  
3 occurs after the number of transfer operations in the measurement period has occurred,  
4 and further comprising means for starting another measurement period to gather transfer  
5 time data after determining paths to disable.

1           21.    The system of claim 19, wherein transfer time data is gathered by path and  
2 transfer size, and wherein the average cumulative transfer time is calculated for each  
3 enabled path and for at least one transfer size.

1           22.    The system of claim 21, wherein the measurement period comprises a  
2 number of transfer operations for all paths for a transfer size, wherein the determination  
3 to disable paths for a transfer size occurs after the number of transfer operations in the  
4 measurement period has occurred, and further comprising means for starting another  
5 measurement period to gather transfer time data for the transfer size after determining  
6 paths to disable for the transfer size.

1           23.    The system of claim 19, wherein the transfer time is measured from the  
2   time the transfer is sent to the device to the time a response is received from the device  
3   indicating that the transfer completed, further comprising means for adding the transfer  
4   time for a transfer transmitted down the path to the cumulative transfer time for the path.

1           24.    The system of claim 19, further comprising:  
2           means for determining a best average transfer time from the average cumulative  
3   transfer times for all paths for each enabled path, wherein the means for determining  
4   whether the average cumulative transfer time for one path satisfies the threshold  
5   comprises means for determining whether the average cumulative transfer time for the  
6   path exceeds the best average transfer time by a percentage amount.

1           25.    The system of claim 24, wherein the means for determining whether the  
2   average cumulative transfer time satisfies the threshold further comprises means for  
3   disabling the path for a first number of transfer operations if the average cumulative  
4   transfer time for the path exceeds the best average transfer time by a first percentage  
5   amount and disabling the path for a second number of transfer operations if the average  
6   cumulative transfer time for the path exceeds the best average transfer time by a second  
7   percentage amount.

1           26.    The system of claim 15, wherein the multiple paths comprise multiple  
2   paths between a first controller and a second controller, and wherein one path is selected  
3   to transmit updates to a primary storage area managed by the first controller to the second  
4   controller to store in a secondary storage area.

1           27.    The system of claim 25, wherein transfer time data is gathered by path and  
2   a size of the update, wherein a path is disabled for a given update size and wherein the  
3   path is capable of being enabled for at least one other update size.

1           28.    The system of claim 15, wherein the paths extend through a network.

1           29.    An information bearing medium for selecting one of multiple data paths to  
2 a device, wherein the information bearing medium includes code capable of causing a  
3 processor to perform:

4           selecting one of multiple paths indicated as enabled to transmit data, wherein a  
5 path is indicated as enabled or disabled;

6           gathering transfer time data for each enabled path capable of being selected; and

7           indicating paths as disabled having transfer time data satisfying a threshold,

8 wherein paths indicated as disabled are not capable of being selected to use to transmit  
9 data.

1           30.    The information bearing medium of claim 29, further capable of causing  
2 the processor to perform:

3           indicating one disabled path as enabled after performing a threshold number of  
4 transfer operations.

1           31.    The information bearing medium of claim 30, further capable of causing  
2 the processor to perform:

3           disabling the path for a first threshold number of transfer operations if the path has  
4 a transfer data time satisfying a first threshold; and

5           disabling the path for a second threshold number of transfer operations if the path  
6 has a transfer data time satisfying a second threshold.

1           32.    The information bearing medium of claim 29, wherein transfer time data is  
2 gathered by path and transfer size, wherein a path is disabled for a given transfer size and  
3 wherein one path disabled for one transfer size is capable of being enabled for at least one  
4 other transfer size.



1           33     The information bearing medium of claim 29, wherein gathering transfer  
2 time data further comprises:  
3           for each enabled path, gathering a cumulative transfer time for all transfer  
4 operations during a measurement period through the path and a cumulative number of the  
5 transfer operations during the measurement period; and  
6           for each enabled path determining the average cumulative transfer time for the  
7 measurement period by dividing the cumulative time by the cumulative number of  
8 transfers, wherein the path is indicated as disabled if the average cumulative transfer time  
9 for the path satisfies the threshold.

1           34.    The information bearing medium of claim 33, wherein the measurement  
2 period comprises a number of transfer operations for all paths, wherein the determination  
3 to disable paths occurs after the number of transfer operations in the measurement period  
4 has occurred, and further causing the processor to perform starting another measurement  
5 period to gather transfer time data after determining paths to disable.

1           35.    The information bearing medium of claim 33, wherein transfer time data is  
2 gathered by path and transfer size, and wherein the average cumulative transfer time is  
3 calculated for each enabled path and for at least one transfer size.

1           36.    The information bearing medium of claim 35, wherein the measurement  
2 period comprises a number of transfer operations for all paths for a transfer size, wherein  
3 the determination to disable paths for a transfer size occurs after the number of transfer  
4 operations in the measurement period has occurred, and further causing the processor to  
5 perform starting another measurement period to gather transfer time data for the transfer  
6 size after determining paths to disable for the transfer size.

1           37.    The information bearing medium of claim 33, wherein the transfer time is  
2    measured from the time the transfer is sent to the device to the time a response is received  
3    from the device indicating that the transfer completed, and further causing the processor  
4    to perform adding the transfer time for a transfer transmitted down the path to the  
5    cumulative transfer time for the path.

1           38.    The information bearing medium of claim 33, and further causing the  
2    processor to perform:  
3           for each enabled path, determining a best average transfer time from the average  
4    cumulative transfer times for all paths, wherein determining whether the average  
5    cumulative transfer time for one path satisfies the threshold comprises determining  
6    whether the average cumulative transfer time for the path exceeds the best average  
7    transfer time by a percentage amount.

1           39.    The information bearing medium of claim 38, wherein determining  
2    whether the average cumulative transfer time satisfies the threshold further comprises  
3    disabling the path for a first number of transfer operations if the average cumulative  
4    transfer time for the path exceeds the best average transfer time by a first percentage  
5    amount and disabling the path for a second number of transfer operations if the average  
6    cumulative transfer time for the path exceeds the best average transfer time by a second  
7    percentage amount.

1           40.    The information bearing medium of claim 29, wherein the multiple paths  
2    comprise multiple paths between a first controller and a second controller, and wherein  
3    one path is selected to transmit updates to a primary storage area managed by the first  
4    controller to the second controller to store in a secondary storage area.

1           42.     The information bearing medium of claim 29, wherein the paths extend  
2 through a network.

[illegible]

# METHOD, SYSTEM, AND PROGRAM FOR SELECTING ONE OF MULTIPLE PATHS TO COMMUNICATE WITH A DEVICE

## ABSTRACT

Disclosed is a method, system, program, and data structure for selecting one of multiple data paths to a device. A selection is made of one of multiple paths indicated as enabled to transmit data. A path is indicated as enabled or disabled. Transfer time data is gathered for each enabled path capable of being selected. Paths having transfer time data satisfying a threshold are indicated as disabled. Paths indicated as disabled are not capable of being selected to use to transmit data.

FIG. 1

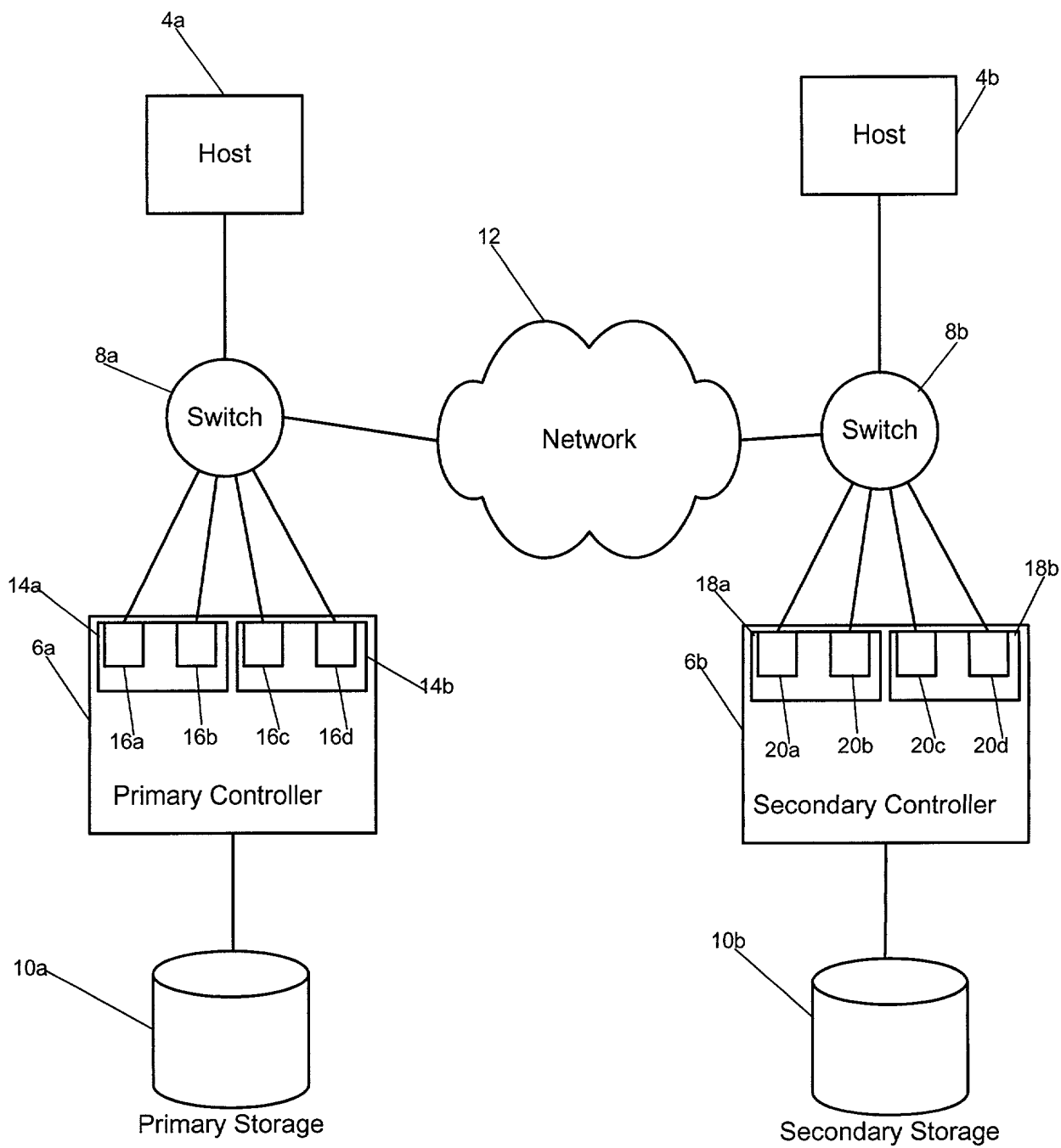


FIG. 2a

Path	Less than 9 Blocks	9-64 Blocks	More than 64 Blocks
1			
2			
3			
...			
...			
16			

FIG. 2b

Block Size Range	Number of Transfers
Less than 9	32
Between 9 and 64	110
Greater than 64	8

FIG. 3

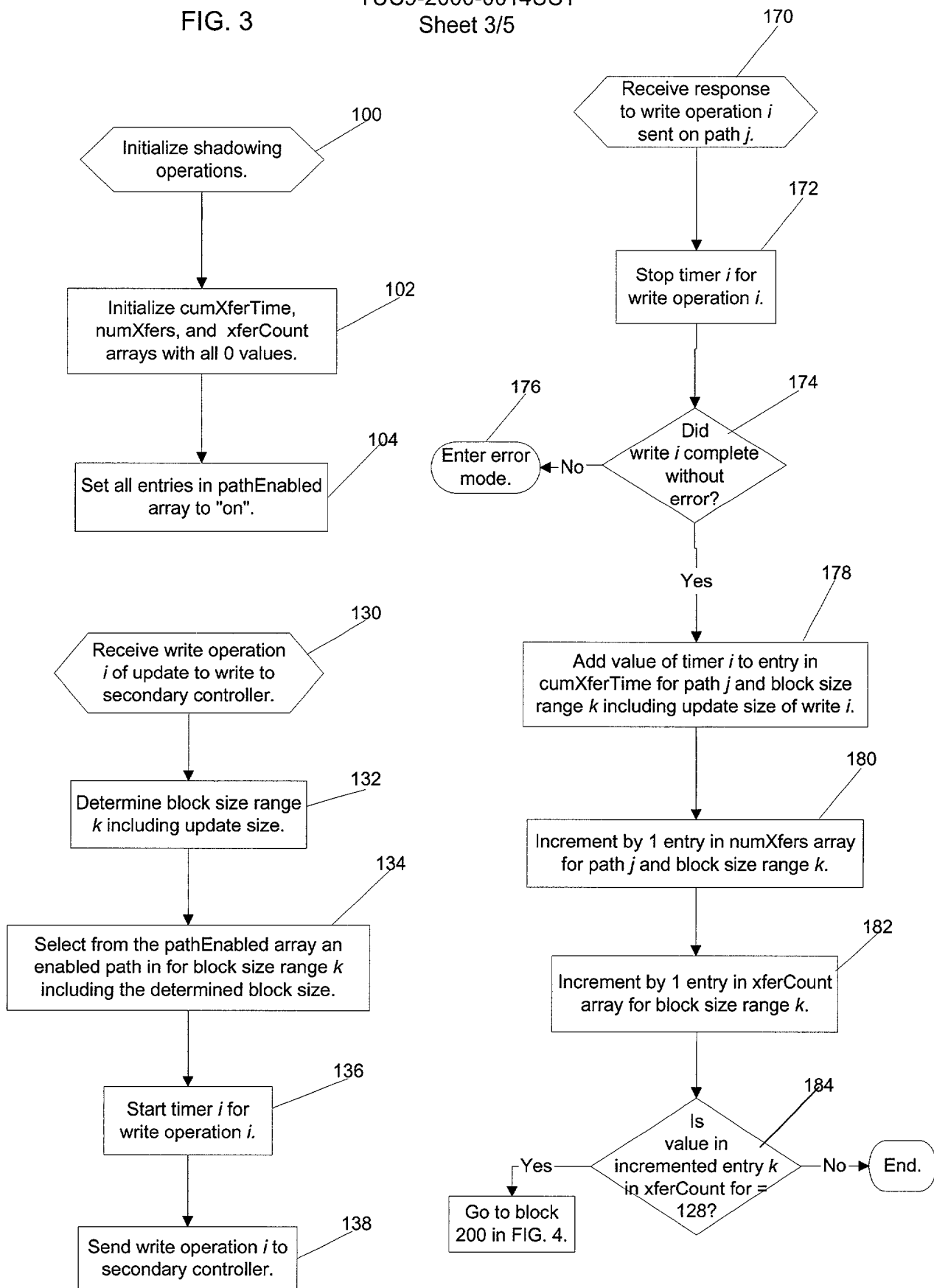


FIG. 4

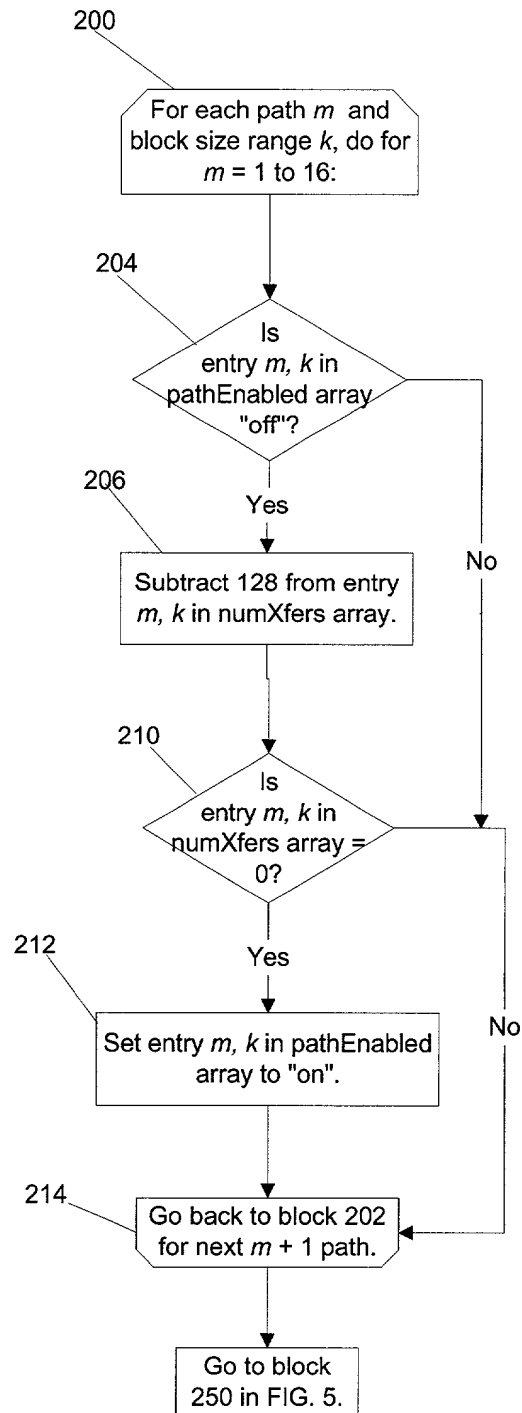
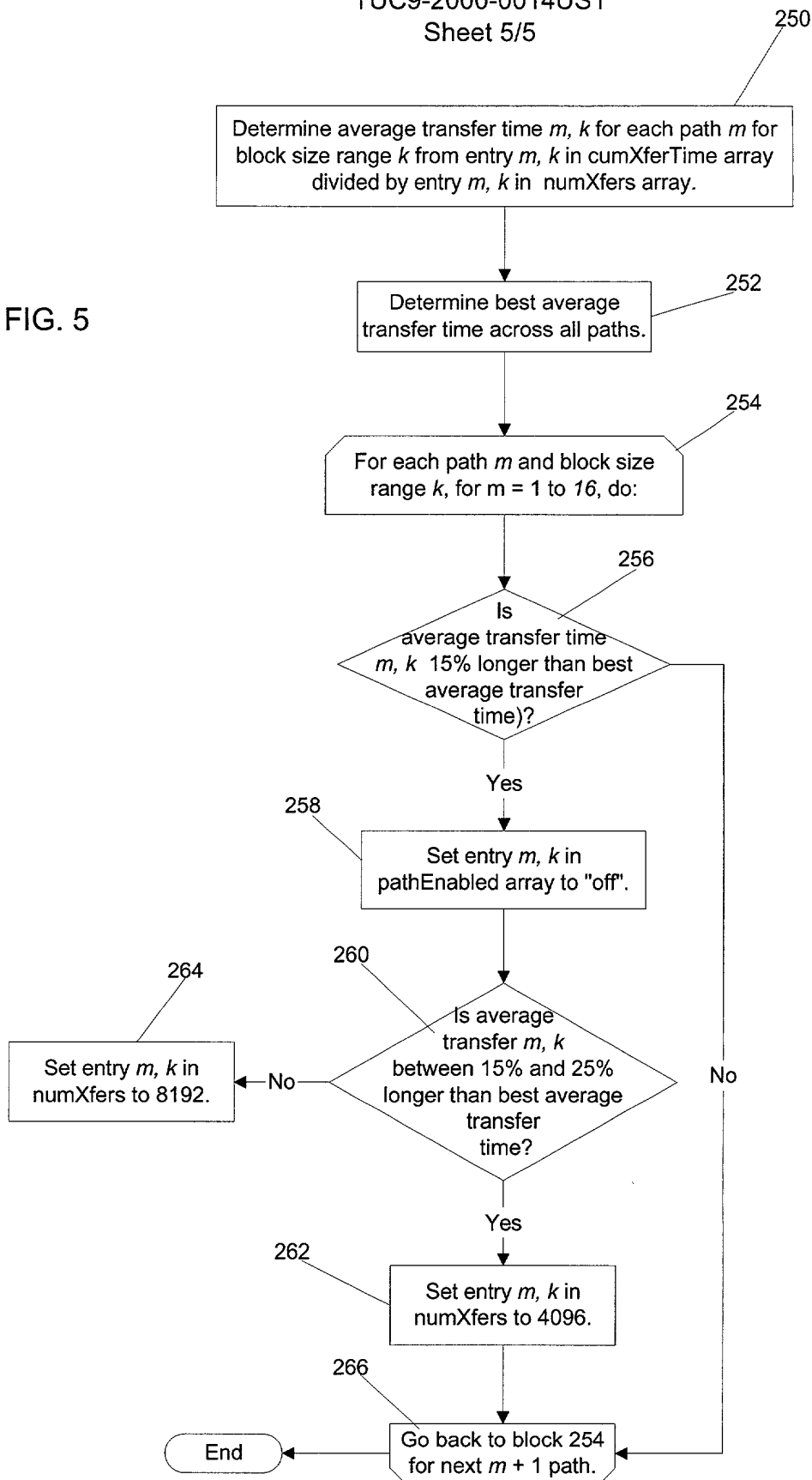




FIG. 5



DECLARATION AND POWER OF ATTORNEY FOR PATENT APPLICATION

DOCKET:  
TUC920000014US1

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name;

I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled

METHOD, SYSTEM, AND PROGRAM FOR SELECTING ONE OF MULTIPLE PATHS TO COMMUNICATE WITH A  
DEVICE

the specification of which (check one)

☒ is attached hereto.

☐ was filed on \_\_\_\_\_

as Application Serial No. \_\_\_\_\_

and was amended on \_\_\_\_\_ (if applicable).

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose information which is material to patentability as defined in Title 37, Code of Federal Regulations, Section 1.56.

I hereby claim foreign priority benefits under Title 35, United States Code, Section 119(a)-(d) or Section 365(b) of any foreign application(s) for patent or inventor's certificate, or Section 365(a) of any PCT International application which designated at least one country other than the United States, listed below and have also identified below any foreign application for patent or inventor's certificate or PCT International application having a filing date before that of the application on which priority is claimed:

Prior Foreign Application(s)

Priority Claimed

None

(Number)

(Country)

Yes No

(Day/Month/Year Filed)

I hereby claim the benefit under Title 35, United States Code, Section 120 of any United States application(s) or Section 365(c) of any PCT International application designating the United States, listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States or PCT International application in the manner provided by the first paragraph of Title 35, United States Code, Section 112, I acknowledge the duty to disclose information which is material to patentability as defined in Title 37, Code of Federal Regulations, Section 1.56, which occurred between the filing date of the prior application and the national or PCT international filing date of this application:

None

(Application Serial No.)

(Filing Date)

(Status) (patented, pending, abandoned)

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

**DECLARATION AND POWER OF ATTORNEY FOR PATENT APPLICATION****DOCKET:**  
**TUC920000014US1**

**POWER OF ATTORNEY:** Joseph F. Villella, Jr., Reg. No. 30,599; David W. Victor, Reg. No. 39,867; William K. Konrad, Reg. No. 28,868; Gary D. Mann, Reg. No. 34,867; Alan S. Raynes, Reg. No. 39,809; Esther E. Klein, Reg. No. 34,337; Robert B. Martin, Reg. No. 26,945; Paik Saber, Reg. No. 37,494; Leslie G. Murray, Reg. No. 31,183; Christopher A. Hughes, Reg. No. 26,914; Joseph C. Redmond, Jr., Reg. No. 18,753; Douglas R. Millett, Reg. No. 31,784; John E. Hoel, Reg. No. 26,279; Edward A. Pennington, Reg. No. 32,588; Robert M. Sullivan, Reg. No. 39,391; G. Marlin Knight, Reg. No. 33,409; Randall J. Bluestone, Reg. No. 40,518; Abdolreza Raissinia, Reg. No. 38,686.

Send correspondence to:

**David Victor, Esq**  
**1180 South Beverly Dr., Ste. 501**  
**Los Angeles, CA 90035**

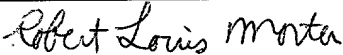
Direct all telephone calls to David Victor at (310) 553-7977

**FULL NAME OF INVENTOR ONE: David Alan Burton**INVENTORS SIGNATURE:  DATE: 5/22/00

RESIDENCE: 7120 South Red Hill Road, Vail, Arizona 85641

CITIZENSHIP: United States of America

POST OFFICE ADDRESS: same as residence

**FULL NAME OF INVENTOR TWO: Robert Louis Morton**INVENTORS SIGNATURE:  DATE: 5/22/00

RESIDENCE: 11560 E. Lusitano Place, Tucson, Arizona 85748

CITIZENSHIP: United States of America

POST OFFICE ADDRESS: same as residence

**FULL NAME OF INVENTOR THREE: Erez Webman**INVENTORS SIGNATURE:  DATE: 6/5/00

RESIDENCE: 14 Degel Reuven Street, Petach-Tikva, Israel 49402

CITIZENSHIP: Israel

POST OFFICE ADDRESS: same as residence